

PATENT ABSTRACTS OF JAPAN

(11)Publication number : 08-328759

(43)Date of publication of application : 13.12.1996

(51)Int.Cl.

G06F 3/06

G06F 3/06

G06F 13/10

(21)Application number : 07-130555

(71)Applicant : MITSUBISHI ELECTRIC CORP

(22)Date of filing : 29.05.1995

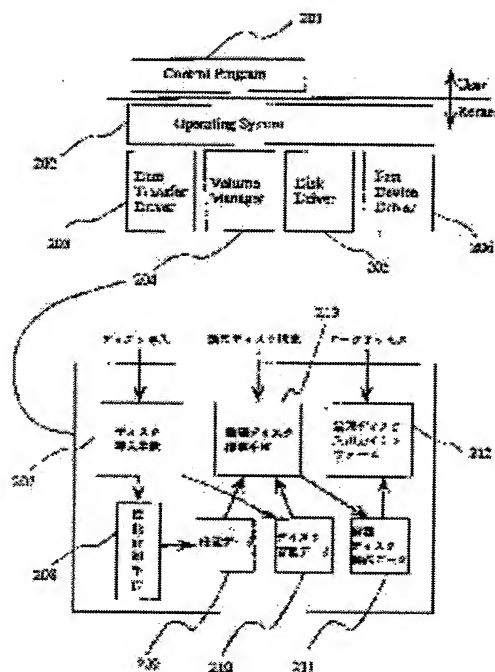
(72)Inventor : SAKAKURA TAKASHI
FUSHIMI SHINYA

(54) INPUT/OUTPUT PROCESSING SYSTEM

(57)Abstract:

PURPOSE: To improve the throughput of the entire system by generating logical disk management data for which a stripe width for uniformizing response time required for the input/output of data for one stripe of respective disk devices is set.

CONSTITUTION: The performance characteristics of the respective disk devices for constituting the stripe are recognized by a performance measurement means 208. It is performed by using a parameter, making the disk devices execute input/output operations and measuring the response time and the result is recorded. Then, based on the measured performance data 209 of the respective disk devices, a data amount to be arranged to the respective disk devices so as to uniformize each response time required for the input/output of the data for one stripe of the disk devices for constituting a disk drive, that is the stripe width, is decided. Then, logical disk constitution data 211 according to it are constituted by a logical disk constitution means 213, that is, the logical disk device of striping constitution is constituted.



LEGAL STATUS

[Date of request for examination] 19.11.1999

[Date of sending the examiner's decision of rejection] 13.02.2001

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number] 3201219

[Date of registration] 22.06.2001

[Number of appeal against examiner's decision of rejection] 2001-03985

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平8-328759

(43) 公開日 平成8年(1996)12月13日

(51) Int.Cl. ⁶	識別記号	庁内整理番号	F I	技術表示箇所
G 0 6 F 3/06	5 4 0		G 0 6 F 3/06	5 4 0
	3 0 1			3 0 1 M
13/10	3 4 0	7368-5E	13/10	3 4 0 B

審査請求 未請求 請求項の数17 O L (全 22 頁)

(21) 出願番号 特願平7-130555

(22) 出願日 平成7年(1995)5月29日

(71) 出願人 000006013

三菱電機株式会社

東京都千代田区丸の内二丁目2番3号

(72) 発明者 坂倉 隆史

鎌倉市大船五丁目1番1号 三菱電機株式
会社情報システム研究所内

(72) 発明者 伏見 信也

鎌倉市大船五丁目1番1号 三菱電機株式
会社情報システム研究所内

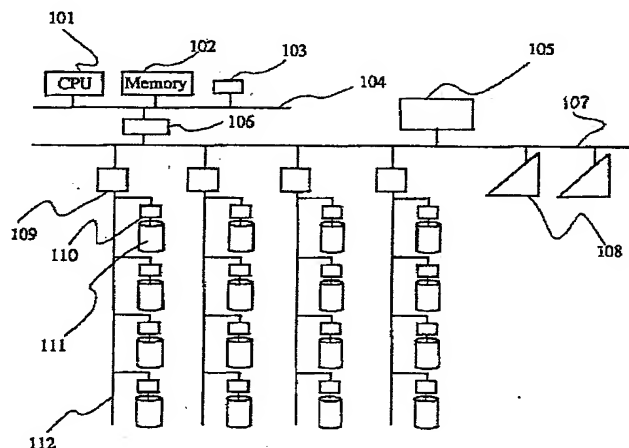
(74) 代理人 弁理士 宮田 金雄 (外3名)

(54) 【発明の名称】 入出力処理システム

(57) 【要約】

【目的】 通常のディスク装置を利用し高速データ転送サブシステムを実現する。

【構成】 ソフトウェア制御によるディスク間同期コストを低減したストライプトディスクとタイムアウトによる同期制御手段、さらに、ユーザメモリを使用しないデータ転送手段により高速大量データ転送システムを実現する。



【 特許請求の範囲】

【 請求項1 】 入出力システムを構成する複数のディスク装置の性能特性をシステム管理者から与えられる性能データもしくは前記ディスク装置を動作させて性能を計測する性能計測手段により収集する性能データ収集手段と、この性能データ収集手段で収集した性能データを基に前記複数のディスク装置を用いて論理ディスク装置を構成する論理ディスク構成手段と、を有する論理ディスク制御手段を備え、

前記論理ディスク構成手段は、前記論理ディスク装置を構成する各ディスク装置の1ストライプ分のデータ入出力に要する応答時間が均等になるストライプ幅を設定した論理ディスク管理データを生成し、前記論理ディスク制御手段は、前記論理ディスク管理データにより前記論理ディスク装置を制御することを特徴とする入出力処理システム。

【 請求項2 】 入出力装置に対する入出力命令発行時に、前記入出力装置の動作に要する制限時間を設定する制限時間設定手段と、設定した制限時間が経過したら前記入出力命令により起動された入出力動作を終了させる手段と、を備えたことを特徴とする入出力処理システム。

【 請求項3 】 前記制限時間設定手段で設定される時間を設定するタイマと、このタイマからの制限時間終了を受け入出力動作を終了させる手段と、を入出力装置の制御をする制御装置に設けたことを特徴とする請求項2に記載の入出力処理システム。

【 請求項4 】 前記入出力装置は複数のディスク装置から構成された論理ディスク装置であることを特徴とする請求項2 または請求項3 に記載の入出力処理システム。

【 請求項5 】 請求項1 に記載の入出力システムにおいて、入出力装置に対する入出力命令発行時に、この入出力命令の処理に関する制限時間を設定する制限時間設定手段と、設定した制限時間が経過したら前記入出力命令により起動された入出力動作を終了させる処理終了手段と、を設けたことを特徴とする入出力処理システム。

【 請求項6 】 前記性能計測手段は、各ディスク装置に対する入出力命令実行時の応答時間が、最短、あるいは、最長となる条件を設定して、入出力命令を実行し、その応答時間を計測することを特徴とする請求項1 または請求項5 に記載の入出力処理システム。

【 請求項7 】 前記性能計測手段は、システムへのディスク装置追加時にこのディスク装置に対する初期化処理の一部としてこのディスク装置の性能測定を行うことを特徴とする請求項1 または請求項5 または請求項6 に記載の入出力処理システム。

【 請求項8 】 前記性能収集手段は、複数のディスク装置が接続された入出力バスの構成、バス転送性能をシステム管理者より与えられたシステム構成データ、あるいは、前記性能計測手段により前記バスに接続された入

力装置を実際に動作させることによってバス性能データを収集し、前記論理ディスク構成手段が、この収集した性能データを基に、各々のディスク性能と接続されたバス転送性能を考慮して、論理ディスク装置を構成することを特徴とする請求項1 または請求項5乃至請求項7のいずれかに記載の入出力処理システム。

【 請求項9 】 前記論理ディスク構成手段は、論理ディスク装置を構成するディスク装置に割り当てられる1回あたりの入出力命令におけるデータ転送量をそのディスク装置の1トラックの中に収めることを特徴とする請求項1 または請求項5乃至請求項9のいずれかに記載の入出力処理システム。

【 請求項10 】 前記論理ディスク構成手段は、外内周部で性能差のある均質な複数のディスク装置により論理ディスクを構成する時は、各々のディスク装置の外内周部を交互に組合せて構成することにより論理ディスクの見かけの入出力性能がデータの配置されたディスク装置上の位置によらず均一となることを特徴とする請求項9に記載の入出力処理システム。

【 請求項11 】 前記論理ディスク制御手段は、論理ディスク装置に対する入出力要求があった場合、その論理ディスク装置を構成する全てのディスク装置への入出力命令発行直前のそのディスク装置のヘッド位置等の状態を動的に判断して、その論理ディスク装置への入出力命令実行時間が最短となるように入出力命令を発行することを特徴とする請求項1 または請求項5乃至請求項10のいずれかに記載の入出力処理システム。

【 請求項12 】 前記論理ディスク制御手段は、複数のディスク装置の各々に複数回の入出力要求が発生する入出力サイズで論理ディスク装置に対して入出力要求があった時は、この論理ディスク装置を構成するディスク装置の性能特性、および、ディスクが接続された入出力バスの特性を判断し、必要に応じて、同一ディスク装置への複数入出力要求をより少ない入出力要求回数とする、あるいは、該複数ディスク装置への入出力要求中に同期ポイントを設けることにより、該論理ディスク装置への入出力時間が最短となるように入出力命令を発行することを特徴とする請求項11に記載の入出力処理システム。

【 請求項13 】 前記論理ディスク制御手段は、システム管理者より与えられた性能を満たす論理ディスク装置を自動的に構成することを特徴とする請求項1 または請求項5乃至請求項12のいずれかに記載の入出力処理システム。

【 請求項14 】 入出力装置に対するデータ転送制御を行うデータ転送ドライバを備え、このデータ転送ドライバは入出力命令の入出力対象装置に論理ディスク装置が含まれている場合は、システムメモリ上にデータ転送用の2面バッファを確保してデータ転送を行うことを特徴とする請求項1乃至請求項13のいずれかに記載の入

3

力処理システム。

【請求項15】 入出力装置の接続される入出力バスの制御をするバスアービタを設け、このバスアービタはデータ転送元及びデータ転送先の装置のIDを登録するソースレジスタ及びデステイニレジスタを備え、前記データ転送ドライバは、入出力命令の入出力対象装置に論理ディスク装置が含まれていない場合は、前記バスアービタを駆動させることによりシステムメモリを利用せずにデータ転送を行うことを特徴とする請求項14に記載の入出力処理システム。

【請求項16】 前記データ転送ドライバは、入出力命令で指定された入力装置と出力装置のデータ転送速度に差がある場合は、システムメモリ上にデータ転送用のバッファを確保して前記バスアービタを駆動することを中心とする請求項14に記載の入出力処理システム。

【請求項17】 前記論理ディスク装置を構成する複数のディスク装置が各々異なる入出力バスに接続されて複数の論理ディスク装置を構成する請求項1または請求項5乃至請求項13のいずれかに記載の入出力処理システムにおいて、前記各入出力バスに接続された各々のディスク装置を制御するディスク制御装置に同じ入出力バスに接続されたディスク装置間での複写手段を設け、異なる論理ディスク装置間でデータの複写を行うことを特徴とする入出力処理システム。

【発明の詳細な説明】

【0001】

【産業上の利用分野】 本発明は計算機システムにおける入出力処理システムに関し、特に大量データを高速に入出力するための入出力処理システムに関する。

【0002】

【従来の技術】 近年、磁気ディスク記憶装置の小型化、低価格化に伴い、RAIDと呼ばれるディスクサブシステムがパーソナルコンピュータやワークステーションに普及してきた。RAIDは複数のディスク装置で論理的なディスク装置を構成して性能や信頼性の向上を図るもので、使用目的によって、一般的にRAID0レベルからRAID5レベルに区分される構成方式のいずれかが採用される。RAIDは一般的にディスクコントローラの機能として実現され、パーソナルコンピュータではSCSIホストアダプタの機能として実現されることが多い。

【0003】 一方、ハードウェアによらずソフトウェアによる制御で、複数のディスク装置により論理的なディスク装置を構成する技術も知られている。IBM社のロジカルボリュームマネージャ、ベリタス社のボリュームマネージャなどが、その代表的な製品であり、これらにおいては、ディスクコンカティネーティング、ディスクミラーリング、ディスクストライピング等の機能が提供されている。また、製品によってはRAID機能を全てカバーすべく試みているものもある。

4

【0004】 RAIDシステムにおいては、大量で高速な入出力のためにはRAID3レベルあるいはRAID0レベルが用いられる。RAID3レベルはバイトインターリーブ構成でパリティデータを付加したデータレイアウト構成をとる。RAID0レベルはソフトウェアによるディスクストライピングと同等でパリティデータが付加されない。

【0005】 ディスクストライピングやRAIDにより期待できる入出力性能上の効果は、インターリーブ配置されたデータに対して複数のディスク装置を含むようなデータサイズの入出力要求があった時に、各々の当該ディスク装置に入出力要求を振り分け、ディスク装置を並列動作させることによってもたらされるユーザから見た入出力性能の向上と、論理ディスク装置に対するコンカレントアクセスを複数のディスク装置に振り分けることによってもたらされるシステムとしてのスループット性能の向上があるが、大量データの高速入出力のためには前者の効果を追求することになる。

【0006】 ところで、ストライプ構成された論理ディスク装置に対する入出力は、それぞれの当該ディスク装置に振り分けられるが、1回の論理ディスク装置に対する入出力ごとに振り分けられた全ての入出力に対して同期をとらねばならない。並列に動作するディスク装置の処理時間にばらつきがあると、処理終了の待ち合わせ時間をそのばらつきの分だけ必要になり、論理ディスク装置の入出力性能の劣化につながる。この同期を乱す要因としては、個々のディスク装置そのものの性能差、シーク（位置決め）、回転待ち等に起因する動的な性能のばらつき、データバスの待ち合わせ等がある。

【0007】 ディスク装置間の待ち合わせのオーバーヘッドを削減するためにRAIDシステムにおいては、例えば、RAID3レベルの適用されたシステムでは、均質なディスク装置でシステムを構成し、ディスクの回転待ちによる同期の乱れをなくすため、全てのディスクの回転を同期させ、各ディスク装置への入出力命令の遅延を吸収するために、データを配置するセクタ位置を少しずつずらすといった技術が知られているが、これには、一つのハードウェアサブシステムとしてRAIDを設計できるという前提が実現条件となる（MICHELLE Y. KIM, "Synchronized Disk Interleaving", IEEE TRANSACTIONS ON COMPUTERS, VOL. C-35, NO. 11, NOVEMBER 1986）。

【0008】 また、さらに細かい視点で、ディスク装置間の同期を乱す要因として、ディスク装置のヘッドの位置ずれ補正をとりあげ、ある基準となるディスクにおいて、位置ずれ補正が必要になった時点で、ディスクアレイを構成する全てのディスク装置に位置ずれ補正を促すコマンドを発行し、位置ずれ補正による同期の乱れを最

5

小限に押さえる技術が特開平5 - 2 7 9 1 0 に開示されている。

【 0 0 0 9 】また、特開平5 - 2 5 7 6 1 1 号公報には、RAI Dシステムにおいて、論理ディスク装置をソフトウェア制御により、柔軟な構成で実現することが開示されている。即ち、ディスクアレイのパーティショニングを目的として、柔軟に1 つあるいは複数のディスクアレイ上にRAI Dレベルの異なる論理ディスク装置を混在させ、なおかつ、ディスク装置間の同期を乱すことによる性能劣化を押さえるデータレイアウト方法について10の技術が開示されている。

【 0 0 1 0 】

【 発明が解決しようとする課題】 以上のように、従来においては、論理ディスク装置をソフトウェア制御により論理ディスク装置を構成するものは存在したが、最適な論理ディスク装置を構成するために、論理ディスク装置を構成するディスク装置の性能を測定して、そのデータを基に、論理ディスク装置を構成する技術はなかった。

【 0 0 1 1 】本発明は、効率のよい入出力処理システムを提供し、システム全体のスループットを向上させることを目的とする。より具体的には、RAI Dを使用せずに、一般的なディスク装置を使用して最適な論理ディスク装置を構成して入出力システムの性能を向上させることを目的とする。また、画像データのような大量データの処理においては、1 回の入出力データ量が非常に大きなものとなるが、この場合1 回の入出力命令に対し指定された大量のデータが完全に転送されるのを待つよりも、指定された量よりも少ないデータを受取り、より早期に処理が開始できる入出力処理システムを提供することを目的とする。さらに、入出力装置間の転送において、効率のよい転送が可能な入出力システムを提供することを目的とする。

【 0 0 1 2 】

【 課題を解決するための手段】 第1 の発明に係わる入出力処理システムは、入出力システムを構成する複数のディスク装置の性能特性をシステム管理者から与えられる性能データもしくは前記ディスク装置を動作させて性能を計測する性能計測手段により収集する性能データ収集手段と、この性能データ収集手段で収集した性能データを基に前記複数のディスク装置を用いて論理ディスク装置を構成する論理ディスク構成手段と、を有する論理ディスク制御手段を設け、前記論理ディスク構成手段は、前記論理ディスク装置を構成する各ディスク装置の1 ストラップ分のデータ入出力に要する応答時間が均等になるストライプ幅を設定した論理ディスク管理データを生成し、前記論理ディスク制御手段がこの論理ディスク管理データにより前記論理ディスク装置を制御するようにしたものである。

【 0 0 1 3 】第2 の発明に係わる入出力処理システムは、入出力装置に対する入出力命令発行時に、入出力装

6

置の動作に要する制限時間を設定する制限時間設定手段と、設定した制限時間が経過したら前記入出力命令により起動された入出力動作を終了させる手段と、をもうけるようにしたものである。

【 0 0 1 4 】第3 の発明に係わる入出力処理システムは、前記制限時間設定手段で設定される時間を設定するタイマと、このタイマからの制限時間終了を受け入出力動作を終了させる手段と、を入出力装置の制御をする制御装置に設けるようにしたものである。

【 0 0 1 5 】第4 の発明に係わる入出力処理システムは、前記入出力装置を複数のディスク装置から構成された論理ディスク装置で構成するようにしたものである。

【 0 0 1 6 】第5 の発明に係わる入出力処理システムは、第1 の発明における入出力システムにおいて、入出力装置に対する入出力命令発行時に、この入出力命令の処理に関する制限時間を設定する制限時間設定手段と、設定した制限時間が経過したら前記入出力命令により起動された入出力動作を終了させる処理終了手段と、を設けようとしたものである。

【 0 0 1 7 】第6 の発明に係わる入出力処理システムは、前記性能計測手段が、各ディスク装置に対する入出力命令実行時の応答時間が、最短、あるいは、最長となる条件を設定して、入出力命令を実行し、その応答時間を計測するようにしたものである。

【 0 0 1 8 】第7 の発明に係わる入出力処理システムは、前記性能計測手段が、システムへのディスク装置追加時にこのディスク装置に対する初期化処理の一部としてこのディスク装置の性能測定を行うようにしたものである。

【 0 0 1 9 】第8 の発明に係わる入出力処理システムは、前記性能収集手段が、複数のディスク装置が接続された入出力バスの構成、バス転送性能をシステム管理者より与えられたシステム構成データ、あるいは、前記性能計測手段により前記バスに接続された入出力装置を実際に動作させることによってバス性能データを収集し、前記論理ディスク構成手段が、この収集した性能データを基に、各々のディスク性能と接続されたバス転送性能を考慮して、論理ディスク装置を構成するようにしたものである。

【 0 0 2 0 】第9 の発明に係わる入出力処理システムは、前記論理ディスク構成手段が、論理ディスク装置を構成するディスク装置に割り当てられる1 回あたりの入出力命令におけるデータ転送量をそのディスク装置の1トラックの中に収めるようにしたものである。

【 0 0 2 1 】第1 0 の発明に係わる入出力処理システムは、前記論理ディスク構成手段が、外内周部で性能差のある均質な複数のディスク装置により論理ディスクを構成する時は、各々のディスク装置の外内周部を交互に組合せて構成することにより論理ディスク装置の見かけの入出力性能がデータの配置されたディスク上の位置によ

50

7

らず均一となるようにしたものである。

【0022】第11の発明に係わる入出力処理システムは、前記論理ディスク制御手段が、論理ディスク装置に対する入出力要求があった場合、その論理ディスク装置を構成する全てのディスク装置への入出力命令発行直前のそのディスク装置のヘッド位置等の状態を動的に判断して、その論理ディスク装置への入出力命令実行時間が最短となるように入出力命令の発行をするようにしたものである。

【0023】第12の発明に係わる入出力処理システムは、前記論理ディスク制御手段が、複数のディスク装置の各々に複数回の入出力要求が発生する入出力サイズで論理ディスク装置に対して入出力要求がある場合は、この論理ディスク装置を構成するディスク装置の性能特性、および、ディスクが接続された入出力バスの特性を判断し、必要に応じて、同一ディスク装置への複数入出力要求をより少ない入出力要求回数とする、あるいは、該複数ディスク装置への入出力要求中に同期ポイントを設けることにより、該論理ディスク装置への入出力時間が最短となるように入出力命令を発行するようにしたものである。

【0024】第13の発明に係わる入出力処理システムは、前記論理ディスク制御手段が、システム管理者より与えられた性能を満たす論理ディスク装置を自動的に構成するようにしたものである。

【0025】第14の発明に係わる入出力処理システムは、入出力装置に対するデータ転送制御を行うデータ転送ドライバを備え、このデータ転送ドライバは入出力命令の入出力対象装置に論理ディスク装置が含まれている場合は、システムメモリ上にデータ転送用の2面バッファを確保してデータ転送を行うようにしたものである。

【0026】第15の発明に係わる入出力処理システムは、入出力装置の接続される入出力バスの制御をするバスアービタを設け、このバスアービタはデータ転送元及びデータ転送先の装置のIDを登録するソースレジスタ及びデスティネレジスタを備え、前記データ転送ドライバは、入出力命令の入出力対象装置に論理ディスク装置が含まれていない場合は、前記バスアービタを駆動させることによりシステムメモリを利用せずにデータ転送を行うようにしたものである。

【0027】第16の発明に係わる入出力処理システムは、前記データ転送ドライバが、入出力命令で指定された入力装置と出力装置のデータ転送速度に差がある場合は、システムメモリ上にデータ転送用のバッファを確保して前記バスアービタを駆動するようにしたものである。

【0028】第17の発明に係わる入出力処理システムは、前記論理ディスク装置を構成する複数のディスク装置が各々異なる入出力バスに接続されて複数の論理ディスク装置を構成する入出力処理システムにおいて、前記

8

各入出力バスに接続された各々のディスク装置を制御するディスク制御装置に同じ入出力バスに接続されたディスク装置間での複写手段を設け、異なる論理ディスク装置間でデータの複写を行うようにしたものである。

【0029】

【作用】第1の発明に係わる入出力処理システムにおいては、性能データ収集手段がシステム管理者が提供する入出力装置の性能データもしくは性能計測手段が実際に入出力データを動作させることにより得る性能データを収集する。得られた性能データを基に論理ディスク構成手段が論理ディスクを構成する各ディスク装置の1ストライプ分のデータ入出力に要する応答時間が均等になるように論理ディスクを構成するための論理ディスク構成データを生成する。こうして得られた論理ディスク構成データを用いて論理ディスク制御手段が論理ディスクを制御する。

【0030】第2の発明に係わる入出力処理システムにおいては、制限時間設定手段により入出力命令出指定させた入出力装置の動作に要する時間を設定し、設定された制限時間が経過すると入出力命令を終了させる手段により入出力動作を未完了状態であっても終了させる。

【0031】第3の発明に係わる入出力処理システムにおいては、入出力装置の制御をする制御装置に設けられたタイマに制限時間を設定し、このタイマから制限時間が経過したことを通知されると同じく制御装置に設けられた入出力動作を終了させる手段により入出力動作をたとえ実行中であっても終了させる。

【0032】第4の発明に係わる入出力処理システムにおいては、複数のディスク装置で構成された論理ディスク装置に対する入出力命令の発行時に制限時間設定手段により入出力動作に要する制限時間を設定し、設定された制限時間を経過すると入出力動作を終了させる手段により動作を終了させる。

【0033】第5の発明に係わる入出力処理システムにおいては、第1の発明により構成された論理ディスク装置に対する入出力命令の発行時に制限時間設定手段により入出力動作に要する制限時間を設定し、設定された制限時間を経過すると入出力動作を終了させる手段により動作を終了させる。

【0034】第6の発明に係わる入出力処理システムにおいては、入出力装置を実際に動作させて性能を計測する性能計測手段が、性能を測定する際に、入出力命令に対する応答時間が最短、あるいは、最長となる条件を設定して入出力命令を実行して応答時間(性能)を計測する。

【0035】第7の発明に係わる入出力処理システムにおいては、入出力装置を実際に動作させて性能を計測する性能計測手段が、システムに初めてディスク装置が導入される時には必ず必要な初期化に合わせて、初期化の一部として性能を計測する。

10

20

30

40

50

【 0 0 3 6 】 第8 の発明に係わる入出力処理システムにおいては、性能収集手段が入出力装置の性能の計測に加えて、入出力装置の接続されたバスの性能をシステム管理者より与えられるバスの性能データもしくは性能計測手段により実際に入出力装置を動作させることにより得る性能データを収集し、論理ディスク構成手段が得られた性能データを基にして論理ディスク構成データを生成する。

【 0 0 3 7 】 第9 の発明に係わる入出力処理システムにおいては、論理ディスク構成手段が論理ディスクを構成する各ディスク装置に対する1 回の入出力命令のデータ転送量が各ディスク装置の1 トラックの中に収まるように論理ディスク構成データを生成する。

【 0 0 3 8 】 第1 0 の発明に係わる入出力処理システムにおいては、論理ディスク構成手段が、性能収集手段で収集した性能データを参照した結果、ディスク装置の内周部と外周部で性能に差があると判断した場合には、各々のディスク装置の外内周部を交互に組み合わせた構成の論理ディスク装置を構成する。

【 0 0 3 9 】 第1 1 の発明に係わる入出力処理システムにおいては、論理ディスク制御手段が、論理ディスク装置をアクセスする場合に、論理ディスクを構成する各ディスク装置のその時点のヘッド位置を認識して、入出力命令としての実行時間が最少となるようにディスク装置への起動順序を決定する。

【 0 0 4 0 】 第1 2 の発明に係わる入出力処理システムにおいては、論理ディスク制御手段が、論理ディスクを構成する各ディスク装置に対する入出力命令が複数回の入出力命令に分割しなくてはならないサイズで要求されたときは、論理ディスク構成データを参照して入出力装置に対する入出力命令数が前記分割数よりも少ない回数で済むように制御する。

【 0 0 4 1 】 第1 3 の発明に係わる入出力処理システムにおいては、論理ディスク構成手段が、論理ディスク構成データを自動的に生成する。

【 0 0 4 2 】 第1 4 の発明に係わる入出力処理システムにおいては、データ転送ドライバが、入出力命令の対象とする入出力装置に論理ディスク装置が含まれている場合には、システムメモリ上に2 面バッファを確保して、アプリケーションプログラムではデータ転送バッファを意識せずに済むデータ転送を行う。

【 0 0 4 3 】 第1 5 の発明に係わる入出力処理システムにおいては、データ転送ドライバが、入出力命令の対象とする入出力装置に論理ディスク装置が含まれていない場合には、データ転送元及びデータ転送先の入出力装置のI Dをバスアービタの対応するレジスタに登録する。以降、バスアービタは、入力側装置と出力側装置とにバスの使用权を与えてシステムメモリを使用せずにデータ転送を行う。

【 0 0 4 4 】 第1 6 の発明に係わる入出力処理システム

においては、データ転送ドライバが、入出力命令の対象とする入出力装置の論理ディスク装置が含まれていないと、且つ入力側装置と出力側装置の転送速度に差があると判断すると、システムメモリ上にデータ転送用のバッファを確保した上で、バスアービタを駆動する。

【 0 0 4 5 】 第1 7 の発明に係わる入出力処理システムにおいては、ディスク制御装置に設けられた複写手段により同一ディスク制御装置に接続されたディスク装置間での複写を行う。

【 0 0 4 6 】

【 実施例 】

実施例1 . 図1 は本発明を適用したハードウェアシステムの構成を示す図である。図において、1 0 1 は本システム全体の処理をするCPU、1 0 2 はCPU1 0 1 が処理する命令及びデータが格納されるメモリ、1 0 3 は時間を計測するシステムタイマで、各々メモリバス1 0 4 に接続されている。メモリバス1 0 4 はブリッジ1 0 6 を介してローカルバス(主たる入出力バス) 1 0 7 に接続されている。このローカルバス1 0 7 はバスアービタ1 0 5 により制御される。また、ローカルバス1 0 7 には、DMAコントローラを備えたSCSI バスアダプタ1 0 9 を介して同じ構成の4 式のSCSI バス1 1 2 と高速デバイス1 0 8 が接続されている。この高速デバイス1 0 8 は、DMAインターフェースを有する装置で、例えば、ソートプロセッサのようにディスク装置から受け取った大量のデータに対して即時的なデータ処理を行い、処理結果をディスク装置へ再び出力するものである。さらに、各SCSI バス1 1 2 には、ディスクコントローラ1 1 0 を介して同じ構成の4 台のディスク装置1 1 1 が接続されている。本実施例では、ディスク装置1 1 1 は同一のものが使用されているが、もちろん、異種のものが接続されても構わない。

【 0 0 4 7 】 図2 は本発明のソフトウェア構成を示す図である。図において、2 0 1 はディスク装置1 1 1、または、論理ディスク装置から、高速デバイス1 0 8 へ、あるいは、逆に高速デバイス1 0 8 からディスク装置1 1 1 または論理ディスク装置へデータ転送を行う制御プログラムである。この制御プログラム2 0 1 は1 ユーザアプリケーションプログラムとして動作する。2 0 2 はオペレーティングシステムであるが、本実施例においては、オペレーティングシステム2 0 2 はフレームワークとして主に利用され、オペレーティングシステム自体の機能はほとんど使用されない。制御プログラム2 0 1 も実質的にはデータ転送ドライバ2 0 3 に制御を依頼する。つまり、「デバイスA のオフセット1 MB 目から3 0 MB のデータをデバイスB に書き込め」という形で指令を出すのみである。データ転送ドライバ2 0 3 は、その名の通り、疑似ドライバとしてカーネルに組み込まれる。2 0 4 はソフトウェア制御によって、複数のディスク装置を使用して論理ディスク装置を構成して、オーバ

11

ヘッドを排して、より小さい応答時間、つまり、より高速な大量データ転送、あるいは、ある制限時間内での確実な応答を得るためのソフトウェアである（このソフトウェアが論理ディスク制御手段の一部で、以後、このソフトウェアをボリュームマネージャと呼ぶ）。205はディスク装置111のコントロールを行うディスクドライバ（論理ディスク制御手段の一部）、206は高速デバイス108のコントロールを行うデバイスドライバ（論理ディスク制御手段の一部）である。なお、バスアービタ105のデータ転送機能を利用する場合のコントロールは、データ転送ドライバ203が直接行う。なお、論理ディスク制御手段は、論理ディスク装置を制御するだけでなく、論理ディスク装置以外の入出力装置も制御する。

【0048】ここで、本発明の中核となるボリュームマネージャ204の機能について説明する。ボリュームマネージャ204は、複数のディスク装置により論理ディスク装置を構成することにより、実際のディスク装置よりも大きな連続データ領域を持つ論理ディスク装置を提供したり、論理ディスク装置への書き込みを複数のディスク装置に反映することによるディスクミラーリングや、ストライプ構成された論理ディスク装置を構成して、入出力性能の高速化等を行う。つまり、ボリュームマネージャ204は自身の管理下におかれたディスク装置からシステム管理者の指示によって論理ディスク装置を構成し、ユーザは論理ディスク装置に対して、アクセスを行い、ボリュームマネージャ204の管理するディスク装置に直接アクセスすることはない。もちろんシステム管理者は全てのディスク装置をボリュームマネージャ204の管理下に置く必要はない。

【0049】ボリュームマネージャ204について、もう少し説明しておく。ボリュームマネージャ204は、ディスク装置をシステムに導入するディスク導入手段207、このディスク導入手段207がディスク導入時にディスク装置から読み取ったディスクラベルとシステムのディスク属性ファイルから生成されるディスク管理データ210、このディスク導入手段が起動するディスク装置の性能を測定する性能計測手段、この性能計測手段で測定した性能データ209及び、この性能データ209とディスク管理データ210とから論理ディスク構成データを生成する論理ディスク構成手段213、論理ディスク構成手段213で生成された論理ディスク構成データを用いて論理ディスク装置へのアクセス時のインタフェースを司る論理ディスク入出力インタフェースから構成される。

【0050】ここで、ディスクストライピングという技術につき整理しておく。理解を容易にするために、2台の全く均等なディスク装置が、十分に高速な入出力バスに接続されているとすると、あるサイズの連続した大量データは、例えば128セクタ分のサイズ（以後、この

12

データ量単位をチャンクと呼ぶ。）に区切られ、チャンクは1番目から交互に、各々のディスク装置に順番に格納される。このようにストライプ構成された論理ディスク装置の先頭に2つのチャンクを含む入出力サイズで入出力要求を発行すると、入出力要求は2つのディスク装置に振り分けられ処理されるので、見かけの入出力性能は2倍になる。

【0051】さて、例えば、システム管理者が通常のディスク装置の2倍の性能のディスク装置が必要となり、通常のディスク装置2台とボリュームマネージャ204を利用してこれを実現する場合で説明する。システム管理者は、ボリュームマネージャ204をシステムに導入する。ボリュームマネージャ204は、ディスク導入手段207に対して、新規2台のディスク装置の導入指示を与え、ボリュームマネージャ204の管理下に登録する。登録要求を受けるとボリュームマネージャ204はディスク装置111の初期化を行い、論理ディスク装置とするための管理メタデータ（ディスク管理データ210）を書き込む。この状態では、まだ、論理ディスク装置としての構成は、なされていない。システム管理者はさらに、ボリュームマネージャ204に2台でストライプ構成の論理ディスク装置を構成するように、論理ディスク構成手段213に依頼する。論理ディスク構成手段213は論理ディスク構成データ210を生成し、これによって初めてユーザは論理ディスク装置を利用できるようになる。ユーザが論理ディスク装置にアクセスする場合は、論理ディスク入出力インタフェース212にオペレーティングシステム202を経由してアクセスする。入出力要求を受けるとボリュームマネージャ204は論理ディスク装置の論理ディスク構成データを参照して、論理ディスク装置を構成する実際のディスク装置111に対し、入出力命令を出す。

【0052】以下、順次本入出力処理システムの動作につき説明していく。最初に、システムへのディスク装置を導入する時に行われる、ディスク装置単体の性能の計測、バス能力の計測等について説明する。先ず、概要について説明する。これは、ディスク装置等の性能を性能計測手段208（性能データ収集手段の一部）によって、ディスク装置の入出力サイズによる応答性能、ディスクアドレスによる応答性能、ディスクキャッシュ効果等の推定を行い、ストライプを構成する各ディスク装置の性能特性を知るものである。具体的には、上記パラメータを用いてディスク装置に対して、入出力動作を実行させてその応答時間（性能データ209の一部）を測定することによって行い、その結果を記録する。そして、その測定した各ディスク装置の性能データ209を基に、ディスクストライプを構成するディスク装置の1ストライプ分のデータ入出力に要する応答時間が各々均等になるように各ディスク装置に配置すべきデータ量、つまり、ストライプ幅を決定し、それに沿った論理ディス

13

ク構成データ211を論理ディスク構成手段213により、即ちストライピング構成の論理ディスク装置を構成する。こうして構成された論理ディスク装置に対しての入出力要求は、ストライプを構成する各々のディスク装置に振り分けられ、各々の入出力装置は均等な時間で応答するので、大まかには、単純にディスク装置に対して入出力する場合に比較して、見かけの応答時間はストライプを構成するディスク装置台数分の1となる。

【0053】ボリュームマネージャ204は、その後にシステム管理者によって要求される論理ディスク装置構成要求時にこの計測された性能データ209を基に、要求された論理ディスク装置性能を確保できるかどうかを判断する。ディスク装置の導入はディスク装置接続の後システムを立ち上げ、ディスク接続確認後、フォーマット（初期化）を行うコマンド投入によって行われるが、書き込みデータとしてフォーマット用データを用いることにより、ディスク性能測定とディスク装置の初期化を兼ねる。これは、単なるフォーマットではなくて、ボリュームマネージャ204の管理下に置くための、メタデータ（ディスク管理データ210）が書き込まれる。図3に示すのはディスク装置111に対して行った性能測定の結果である。これは前記メタデータとして該ディスク装置上に格納されたものをそのままきり出したスナップショットである。

【0054】図3において、501はシステムの中で該ディスク装置111を特定する名前、502は、このディスク装置111の型式を特定するインデックス情報である。ボリュームマネージャ204はディスク装置111の諸元についてのカタログデータを持ったデータベースをシステム管理者から与えられて、あるいは、ディスク装置111のラベル情報を読み出すことにより、この情報を得る。これにより、ディスク回転数、各シリンダにおける、1トラックあたりのセクタ数などを知ることができる。503はディスク装置111の容量、504はシーク、回転待ちなしの場合のアクセスタイム、505は最も遠い位置からシークし、1回転待った場合のアクセスタイムで、マイクロ秒で記されている。この計測は、基準位置へコマンド発行した後、次セクタ、あるいは、最も遠いセクタへのアクセス命令を発行し、基準位置へのアクセス完了から、計測位置へのアクセス完了までの経過時間を測定することによって行う。506から507までの10個のデータは0シリンダから、100MB毎の位置で転送速度を計測したものがKB/sで記されている。これらのディスク性能データ209は後に説明する論理ディスク構成の際に参照/利用される。なお、本実施例で使用したディスク装置は読み出し性能と書き込み性能に差がないものとする。

【0055】一般に、システムにディスク装置を新規に導入する時には、不良セクタの発見や、ラベル情報の書き込みのために、フォーマットと呼ばれる処理を行う

14

が、この実施例ではこのフォーマット時に、上記のように性能計測手段208によりディスク装置の性能データ209を収集する作業を兼ねて行うことができるにしているので、性能計測に要するコストを低減でき、システム管理者の負担を低減できる。

【0056】次に、入出力バスのスループット性能の測定について説明する。この測定は、ディスク装置の入出力性能や同期動作を妨げる要因となり得る入出力バスのデータ転送性能（性能データ209の一部）や同一バス上での競合を考慮して、論理ディスクの構成を行うためのものである。システム管理者により与えられた、入出力バスのデータ転送能力及び構成データから、もしくは、複数の入出力バスに各々接続された全てのディスク装置を含めた入出力装置を組み合わせさせて駆動し、負荷をかけ、各々の入出力装置の入出力動作の応答時間を測定することにより、各々の入出力バス、さらにそれらが接続された上位の入出力バスのデータ転送能力に余裕があるか、あるいは、転送能力の限界値はどれほどかを知り、ストライプ構成の論理ディスク装置への入出力時に、並列動作するディスク装置のデータ転送により共有する入出力バスのデータ転送能力の限界となることを回避するように、論理ディスク装置をストライプ構成するディスク装置を選択して構成することにより、論理ディスク装置への入出力時に入出力バスネックとなることによる性能劣化を回避するためのものである。

【0057】従って、この測定は、単なるスループット性能だけでなく、入出力バス競合の発生条件を測定対象バスに接続された入出力装置に対して、起動する順番によって規定して、実際にバス競合を発生させ、各々の入出力装置の入出力レイテンシを測定するものである。これにより、バスアービトレーションポリシー等を知り、ストライプを構成するディスク装置の入出力起動時に不要なバス競合を避けるための、起動遅延を挿入するといった適用が可能であるが、本実施例では、ディスク装置が接続されるSCSIバス112、ローカルバス107のスループット性能を計測して、ストライプ構成される論理ディスク装置に対する入出力時に、各入出力バスのスループット性能の限界を超えないようにストライプを構成する。SCSIバス112については、接続された入出力装置、本実施例では、たまたま、同一のディスク装置4台が全てのSCSIバス112に同様に接続されているが、該バスに接続された、入出力装置を全て起動して一定時間に得た入出力データの総量を計算することにより、スループット性能の限界値を知る。これを、ディスク装置一つを動作させた場合のスループット性能と比較して、何台のディスク装置まで、バス転送速度の制約を受けずに動作可能かを知る。本実施例の場合ディスク装置の最大転送速度が2MB/sで、SCSIバス112の最大転送速度は8MB/sであるのでSCSIバス112に接続された4台のディスク装置を同時に動作

させても、単純にはバス転送速度の制約を受けないことがわかる。

【0058】同様に、ローカルバス107の転送速度を全てのディスク装置、及び、高速デバイス108を動作させることにより計測したところ、40 MB/sで、全ての入出力デバイスの転送速度の合計に相当し、結局、本実施例の入出力バスと入出力装置の構成においては、全ての入出力装置に対してどのように入出力動作を実行しても入出力バスの転送速度による制約を受けないことがわかる。

【0059】ところが、本実施例とは異なり、例えば、SCSIバスの転送速度が本実施例のSCSIバスの転送速度より遅く、そのSCSIバスの転送速度では2台分のディスク装置しか、ディスク装置の最大転送速度では動作できないのであれば、ストライプ構成で同一の論理ディスク装置を構成するディスク装置を3台以上、同一SCSIバス上に接続されたディスク装置では構成しないようにすることになる。

【0060】次に、論理ディスク装置の構成方法について説明する。図4は、ストライプ構成の論理ディスク装置を構成する論理ディスク構成手段213の動作を説明するフローチャートである。以下、この図を参照しながら説明する。ステップ601において、ユーザ、つまり、入出力装置を管理するシステム管理者は、必要とするストライプ構成の論理ディスク装置の性能、例えば、8 MB/sと、アプリケーションプログラムが論理ディスク装置のアクセスの際に発行する入出力サイズ、例えば、512 KB、及び、必要とする論理ディスク装置のサイズ、例えば、1 GBを入力する。これにより、本入出力処理システムは、指定された仕様を満たす論理ディスク装置の構成を試みる。ここで、入出力サイズを指定するのは、ストライプ数、及び、ストライプ幅が入出力サイズによって制限されるためである。つまり、入出力サイズは、上述のチャンクサイズの倍数でないと、ディスクストライピングによる性能向上が得られないためである。

【0061】次に、ステップ602において、性能計測手段208によって得たディスク装置の性能データ209から、指定された論理ディスク装置の性能を達成するのに何台のディスク装置の並列動作が必要か、つまり、ストライプのWAY数を決定する。WAY数が決定したら、ステップ603で本入出力処理システムの管理するディスク装置上の空き領域を探し、必要なWAY数分の異なるディスク装置上に、ストライプを構成可能かどうか調べる。ただし、指定された性能基準を満たすのに、ディスクストライピングの必要のない場合もあるので、この時は、単にディスク装置上の空き領域を探すことになる。その後ステップ604で、与えられた入出力サイズによって、この入出力サイズをチャンクサイズとするストライプ構成の論理ディスク装置を構成した場合に、

性能上の効果があるか、あるいは、要求された性能水準はディスクストライピングを必要としないのかを判断し、ストライプ幅の制約により効果がない、あるいは、ディスクストライピングの必要のない場合は、ディスク領域を確保の上期待できる論理ディスク装置性能を、ステップ610でユーザに報告して終了する。

【0062】2WAY以上のディスクストライピングが有効な場合はさらに、ステップ605で同種のディスク装置でストライプ構成が可能か否かを判断し、構成することが可能な場合には、ステップ606でこのディスク装置が内外周部で性能差があるか否かを判断し、内外周部で性能差がある場合には、ステップ607でさらにストライプを構成するディスク装置が同一のオフセット位置から使用可能か否かを判断する。そして、同一のオフセット位置から使用可能であれば、ステップ608でディスク装置を交互に逆方向にストライプ構成するようにする。また、ステップ605またはステップ606またはステップ607で否と判断したときは、ステップ608でストライプ幅で個々のディスク装置の応答速度を平均化するようにストライプ構成をする。

【0063】ステップ609における動作を以下に補足して説明する。ディスクストライピングは次のように構成される。まず、チャンクサイズを決定するが、指定された入出力サイズが、論理ディスク装置の性能を達成するのに必要なストライプ幅の総計より大きくない限り、入出力サイズがそのままチャンクサイズとする。一方、入出力サイズが大きすぎる場合は、必要に応じて、入出力サイズの自然数分の1をチャンクサイズとする。いま、入出力サイズとして512 KBが指定されているとする。この場合、例えば、ディスク装置の性能データ209から、平均32 KBストライプ幅の4WAYで性能達成が可能と考えたとチャンクサイズは、512 KBの1/4の128 KBとする。すでに取得しているディスク装置の性能データ209から、ストライプを構成するディスク装置の当該位置のデータ転送速度を得て、その、逆数の比で、当該チャンクを分割して、各々のストライプ幅とする。これにより、該論理ディスク装置に対する入出力動作の実行時には、各々のストライプを構成するディスク装置の応答時間は平均化される。

【0064】次に、ステップ608における動作を補足して説明する。ステップ608におけるディスクストライピングは次のように構成される。この場合も最初にチャンクサイズ、及び、各ディスク装置のストライプ幅が決定されるが、偶数個用意されるストライプ構成されたディスク装置の同一の使用領域の先頭からと、最後尾から交互にストライプが構成され、そのストライプ幅は1トラック分とする。ディスク装置の性能情報の502により示される当該位置のトラックの含むセクタ数より、該先頭位置と最後尾位置の1トラックのサイズを知り、(先頭トラックのサイズ+最後尾トラックのサイズ)×

ストライプWAY数/2 がチャンクサイズとなる。従って、ストライプ幅は各々の入出力時点でトラックあたりのサイズになる。

【0065】以上のように、論理ディスク装置を構成するための論理ディスク構成データ211を生成登録したならば、推奨する入出力サイズ、つまり、内部で決定されたチャンクサイズと期待できる、この論理ディスク装置の性能をステップ610で報告して終了する。あらかじめ与えられた性能水準と入出力サイズにより、ディスクストライピングの必要がない場合、あるいは、効果が

10 ない場合は、推奨する入出力サイズとして、単体ディスク装置に対して、与えられた性能水準を満たすのに必要な入出力サイズを報告し、効果がなければ、期待できる性能として単体ディスクの性能を報告する。

【0066】図5、図6は以上のようにして構成された論理ディスク装置、及び、論理ディスク装置を構成する各ストライプ構成されたディスク装置の管理情報の一部である。これら管理情報は本入出力処理システムの管理する全てのディスク装置111上にメタデータとしてそのコピーが格納され、必要に応じて、本入出力処理シ

20 テムの管理するメモリ102上に読み出される。便宜のため図には一般的な表形式で書かれているが、実際のディスク装置、メモリ上には、プログラムで管理されたテンプレートに従い格納される。

【0067】以下、図5、図6に示す、ステップ608の実行により構成された論理ディスク装置及びストライプ構成されたディスク装置の管理情報について説明する。図5において、701は該論理ディスク装置のデバイスノード名である。アプリケーションプログラムはこのノードをディスク装置としてアクセスする。702は

30 該論理ディスクのサイズ、つまり、容量を示している。703は、該論理ディスク装置の構成情報として4WAYのストライプ構成であることを示し、704は、ストライプを構成するディスク装置とそのディスク装置の領域を示して、各々、c0t0d0-0、c1t0d0-0、c2t0d0-0、c3t0d0-0という名称で特定されている。理解を容易にするために、c0t0d0がディスク装置を特定し、-0が領域番号を特定する

ようにしているが、これは論理的なものである。ハードウェア上のアドレスを特定するものではない。また、図6には、704に示すディスク装置の領域情報、つまり、リージョン情報へのインデックスc0t0d0-0によって特定されるリージョン情報が示されている。801は該リージョン名、802は該リージョンが配置される物理ディスク装置名、803はリージョンの開始オフセット、804は該リージョンのサイズ、805は該リージョンにおけるデータの配置方向、805には該リージョンにおけるストライプ幅を示している。

【0068】次に、図5、図6に示される管理情報(論理ディスク構成データ211)を持つ論理ディスク装置

に対して入出力動作を要求された時の動作につき説明する。アプリケーションプログラム(制御プログラム201もアプリケーションプログラムの1つである)から本論理ディスク装置に対してアクセスする時は、入出力サイズとしてチャンクサイズ、及び、論理ディスク装置に対するオフセット値が与えられる。言い替えると、アプリケーションプログラムはあらかじめ、ボリュームマネージャ204に対してチャンクサイズを問い合わせることができる。ボリュームマネージャ204は問い合わせ要求を受けると指定された論理ディスク装置のチャンクサイズ等の属性を報告する。本入出力処理システムは、入出力要求を受けると、論理ディスク装置の管理データ211を参照して、この論理ディスク装置は4WAYのストライプ構成で、704で示されるリージョンから構成されていることを知る。また、リージョン情報からこの論理ディスク装置がストライプ構成されていることを知ると、入力された論理ディスク装置に対するオフセット値から、このオフセットはストライプ何チャンク分に当たるかを計算し、ストライプを構成する各々のディスク装置に対し、配置方向がリバースである、つまり、逆方向のリージョンについては、該リージョンの最後尾からオフセットチャンク数分のトラックサイズの累計をオフセットとし、該トラックのトラックサイズを入出力サイズとして、該ディスク装置に入出力命令を発行する。順方向のリージョンに対してはリージョンの先頭から、同様に、オフセットを計算し、該トラックのトラックサイズを入出力サイズとして入出力命令を発行する。論理ディスク装置を構成する全てのリージョンについて、ディスクドライブ205の保持する該ディスク装置のヘッド位置から、最もヘッド位置の遠いものから入出力を起動し、その全ての入出力の完了を待つことによって、この論理ディスク装置への入出力動作は終了する。ただし、リージョン内の各々の位置のトラックサイズはすでに説明したように、本入出力処理システムがディスク装置の種類毎に備えるデータベースによって得る。

【0069】以上のように、この実施例によれば、システムに接続されているディスク装置の性能特性を性能収集手段がシステム管理者の与えるデータもしくは性能計測手段により収集し、収集したデータを基に、論理ディスク構成手段により論理ディスクを構成し、論理ディスク制御手段により論理ディスクを動作させるので、論理ディスクを構成する各ディスクの1ストライプ分のデータ入出力に要する応答時間が均等となり、論理ディスク装置の性能向上を図ることができる。

【0070】また、システムにディスク装置を新規に導入する時には、不良セクタの発見や、ラベル情報の書き込みのために、必ず実行されるフォーマット処理を兼ねて、性能計測手段208によりディスク装置の性能データ209を収集する作業を行うことができるにしている。性能データ209の収集に要するコストを低減す

ることができると共にシステム管理者の負担を低減できる。

【0071】また、実際にシステムに装着されているディスク装置各々の、最大アクセス時間、最小アクセス時間、及び、転送速度を、実際に個々の装置に対して入出力動作を実行し、その応答時間を測定して、測定したデータを一覧表にしたテーブルを用いて論理ディスクを構成するので最適な論理ディスク装置を構成できる。

【0072】さらに、ディスク装置の入出力性能や同期動作を妨げる要因となり得る入出力バスのデータ転送性能や同一バス上での競合を考慮して、論理ディスク装置を構成するために、システム管理者により与えられた入出力バスのデータ転送能力及び構成データを用いて、もしくは、性能計測手段により複数の入出力バスに接続された全てのディスク装置を含めた入出力装置を組み合わせ動作させることにより、負荷をかけ、各々の入出力装置の入出力動作における応答時間を測定することにより、各々の入出力バス、さらにそれらが接続された上位の入出力バスのデータ転送能力に余裕があるか、あるいは、転送能力の限界値はどれほどかを知り、ストライプ構成の論理ディスク装置への入出力時に、並列動作するディスク装置のデータ転送により共有する入出力バスのデータ転送能力の限界となることを回避するように、論理ディスク装置をストライプ構成するディスク装置を選択して構成するようにしたので、構成された論理ディスク装置に対する入出力動作時に入出力バスネックとなることによる性能劣化を回避することができる。

【0073】また、論理ディスク装置のストライプを構成する1チャンク分のデータをディスク装置上の1トラックに収まるように配置することにより、ヘッドシークのコスト（位置決めに要する時間）をシーケンシャルアクセスにおいて、その方向によらず平均化する手段を設けるようにしたので、外周から内周方向へシークしても、内周から外周方向にシークしても、読みだし時のディスクキャッシュ効果を除き差がなくなるので、より柔軟にディスクストライプを構成できる。

【0074】さらに、内外周部で性能差のある複数のディスク装置でストライプ構成する場合には、内外周部が交互に組み合わせられるようにストライプを構成することにより、ストライプ構成された論理ディスク装置の読みだし／書き込み位置に関わらず、均一な入出力応答時間を得ることができる。

【0075】さらに、ストライプ構成された論理ディスク装置に対して入出力要求があった場合に、そのアクセス要求位置へのディスクヘッドの現位置からの距離と各ストライプを構成するディスク装置の性能的に有利な入出力サイズを考慮して、その論理ディスク装置としての入出力動作が最も短時間に終了するように、全ての必要とされるディスク装置への入出力命令の発行順序を制御するようにしたので、論理ディスク装置としての入出力

動作の効率が向上する。

【0076】さらに、論理ディスク装置を構成するWAY数、ストライプ幅等の要素を自動的に生成するようにしたので、システム管理者の負荷を軽減し、容易に最適な論理ディスクを構成できる。

【0077】実施例2. 本実施例は、実施例1で説明したようにして構成された論理ディスク装置に対して、複数のチャンクを含む入出力サイズで入出力要求があった場合、同一のディスク装置に対する複数の入出力動作を取りまとめて行うものである。この実施例は、ストライプ構成された論理ディスク装置に対して複数のチャンクを含む入出力サイズで入出力要求があった時、同一ディスク装置に対して発行される複数回の入出力命令をを1回にまとめた方が、その論理ディスク装置への入出力動作がより速く終了するならば、全体性能として、論理ディスク装置の入出力性能を向上させることができることを前提にするものである。しかしながら、論理ディスク装置を構成するディスク装置、及び、その接続された入出力バスの優先度決定方式、あるいは、ディスク装置をドライブするデバイスドライバの特性によって、必ずしも、ディスク装置への入出力命令をとりまとめて発行することが、全体性能の向上につながらないことがあり、このような場合においては、必要に応じて積極的にディスク装置間の動作を同期する、同期ポイントを挿入することにより、論理ディスク装置に対する入出力時間の向上を図る。

【0078】本実施例においては、複数の入出力命令を纏めて1つの入出力命令とするので、所謂データチェーニング機能を必要とするが、ディスクドライバ205、及び、SCSIバスアダプタ109に、メモリとのデータ転送時のスキッピング／ギャザリング機能を持たせ、DMAリストにより、このデータチェーニングを実現している。

【0079】実施例1で説明した論理ディスク装置では、シーク方向が逆になるデータ配置方向のリージョンを含んでいるので、必ずしも入出力命令ををまとめることが性能的に有利にならないため、この入出力命令の取りまとめは行わない。もし、実施例1の論理ディスク装置が図4で説明したステップ609の実行により構成された論理ディスク装置であれば、本実施例において、入出力バスは転送能力に十分余裕を持つので、有利であると判断され、同一ディスク装置上の複数の入出力命令の取りまとめが行われる。

【0080】実施例3. 本実施例は、本発明の入出力処理システムにおける高速処理を実現するためのさらに他の実施例である。図7はその実現手段をディスクコントローラ110上に設けた構成を示すブロック図である。図において、301はSCSI制御部、304はディスク制御部、305はタイマ部である。本実施例においては、ディスクコントローラ110上にタイマ部30

5 を設けるようにしたが、メモリバス104に接続されたシステムタイマ103を用いることにより実現することも可能である。

【0081】まず、本実施例におけるディスクドライブ205とディスクコントローラ110の動作について説明する。ディスクコントローラ110はSCSIベンダユニークコマンド(制限時間設定手段の一部)によりディスクドライブ205から、SCSIバスアダプタ109、SCSIインターフェース301、SCSI制御部302を経由して、タイマ部305(制限時間設定手段の一部)にセットするタイムアウト値を受け取ることができる。該コマンドによってタイマ中のカウンタレジスタ(図示せず)にタイムアウト値がセットされた後、リード、または、ライトコマンドをイニシエータ(コマンドの発行元)から受け取れることを契機にカウントダウンを始める。カウンタレジスタが0になるとタイマ部105はカウントダウンを停止し、SCSI制御部302に対してカウント値が0になったことを報告する。この報告を受けるとSCSI制御部302はコントローラ内部バスのアボートラインをアサートする。また、ディスク制御部304はいかなるフェーズであっても、SCSI制御部302から参照可能な内部レジスタ(図示せず)に、中断時フェーズ、転送済みデータ量などを記録して動作を中止する。SCSI制御部302は、SCSIバス112上でディスコネクトされていればイニシエータとリコネクトの上ステータスフェーズに遷移して、転送済みデータ量(ステータス情報)をメッセージとして送出する。なお、全てのデータ転送完了時になおタイマ動作中であれば、SCSI制御部302がカウンタレジスタを0にリセットする。

【0082】次に、本実施例の詳細な動作を説明する。ディスクドライブ205はタイムアウト値を設定時に次のように動作する。ディスクドライブ205の使用者プログラム、即ち本実施例におけるデータ転送ドライブ203、あるいは、ボリュームマネージャ204は、ディスクドライブ205にデータ転送時のタイムアウト動作を要求するときには、ディスクドライブ205のデータ転送サブルーチンコール時に、ディスクドライブ205が自由に利用できる、システムに共通なブロックデバイスコントロールテーブル(ディスク装置のようにブロック単位で転送を行うデバイスを制御するためのテーブル)のフィールドを使用してタイムアウト値をセットする。タイムアウト値がセットされるとディスクドライブ205はSCSIバスアダプタ109に対し、リード/ライトコマンドの発行を要求する前に、該タイムアウト値を渡し、上述のベンダユニークコマンドの発行を要求する。続いて、リード/ライトコマンドの発行を要求するとディスクドライブ205は、SCSIバスアダプタ109から処理の終了を割り込みによって通知される迄スリープ状態となる。処理終了の通知を受けると、ディ

スクドライブ205は、上述のコントロールテーブル中のタイムアウト値を0にリセットし、転送済みのデータ転送量をコントロールテーブル中に用意されている、データ転送量(ステータス情報)として報告する。

【0083】次に、本実施例をストライプ構成された論理ディスク装置に対して適用した具体的な例について説明する。ところで、本実施例における入出力処理システムは制御プログラム201によりディスク装置をストライプ構成として高速化し、高速デバイス108へのデータ供給、また、高速デバイス108からの出力を論理ディスクディスクへの書き込みを行うが、高速デバイス108の処理能力を有効に活用するために、特にデータ転送開始時に当たっては、実際に送られるデータ量よりも、より、処理開始までのレイテンシを小さくすることに主眼をおくものである。また、もともと大量データを扱うので、多少入出力オペレーションの回数が増えても全体性能に影響しない。

【0084】制御プログラム201がデータ転送ドライブ203を通して、必要と判断すると、ストライプ構成の論理ディスク装置に対して、タイムアウト値と入出力サイズを与える。ボリュームマネージャ204はタイムアウト値を上述のブロックデバイスコントロールテーブルを通して受け取ると、ストライプを構成するディスク装置の入出力起動時に各々のディスク装置について用意されるコントロールテーブルのタイムアウトフィールドにデータ転送ドライブ203から受け取ったタイムアウト値を設定する。図8は4台のディスク装置でストライプ構成された論理ディスク装置に対し、タイムアウト値を設定して読み出しを実行したときに、タイムアウト発生後、ボリュームマネージャ204が4台のディスク装置に対して同期を取り、各々ディスク装置からデータ転送量を受け取った状態を模式的に表したものである。図8において、901は第1のディスク装置、902、903、904はそれぞれ、第2、第3、第4のディスク装置を表し、上下方向は各ディスクのデータ転送の進捗を表す。横方向の点線は該論理ディスクの1チャンク分のデータ量を表す。従って、この例においては、ストライプ構成された論理ディスク装置に対して4チャンク分が入出力サイズとして渡されたことになる。本発明による入出力処理システムにおいては、実施例2で、すでに説明したように、可能で、かつ、有効ならば複数の同一ディスク装置への入出力命令は取りまとめられる。この図8において、各ディスク装置には4チャンク分の入出力命令が取りまとめられて渡され、この、入出力命令に対してタイムアウト値が渡される。既に、説明したように、ディスクドライブ205、及び、SCSIバスアダプタ109には、データチェーニング機能を有しているので転送先または転送元のメモリアドレスが連続していなくても問題はない(図8において平面をメモリ空間と見立てると横方向が連続である)。

【 0 0 8 5 】 タイムアウト 発生後、ボリュームマネージャ 2 0 4 が 4 台のディスク装置 9 0 1 ~ 9 0 3 の同期をとった後の図 8 に示される各ディスク装置の転送状況は、9 0 4 に示す第 4 のディスク装置を除いて未完である。しかし、いずれのディスク装置も 9 0 5 で示す第 3 チャンクまでについてはデータ転送を完了しているのでボリュームマネージャ 2 0 4 は 3 チャンク分のデータ転送が完了したことを、データ転送ドライバ 2 0 3 に報告する。本実施例においては、論理ディスク装置に対して、余りに短いタイムアウト 設定をするとデータ転送量が 0 となる可能性があることがわかる。従って利用者としては、入出力サイズとして、入出力要求を出す論理ディスク装置の複数チャンク分を与え、それに見合ったタイムアウト 値を与えることが有効である。本実施例では、制限時間設定による入出力動作の終了を実施しているが、図 8 に示すようにほぼ全てのディスク装置のデータ転送は同期している。

【 0 0 8 6 】 複数チャンク分の入出力サイズでリード/ライトを行った時にストライプを構成する同一のディスク装置に対する入出力の取りまとめが、その、有効性から行われない場合には、1 チャンク分ずつ、ボリュームマネージャ 2 0 4 はディスク装置間の同期をとる。全てのディスク装置への入出力命令発行時には上位指定のタイムアウト 値を設定するが、第 2 チャンク以降の入出力命令発行時には、第 1 チャンクの入出力命令発行から同期取得時点までの経過時間をシステムタイマにより計測し、これを該タイムアウト 値から減じて、第 2 チャンクの入出力命令発行時にタイムアウト 値として設定する。以降、該タイムアウト 値が 0 になるまでこれを繰り返す。なお、本実施例においては、SCSI のベンダユニークコマンドを用いて時間設定をする例を示したが、必ずしも、ベンダユニークコマンドを使う 必要はなく、また SCSI インタフェースを用いる 必要もない。

【 0 0 8 7 】 以上のように、入出力動作の制限時間を設定する制限時間設定手段と 設定した制限時間が経過すると入出力動作を終了させる手段を設け、指定されたデータ量の入出力動作の完了前に、指定された制限時間に達した時は、入出力動作を終了させ、その時点での転送済みデータ量を報告するようにしたので、アプリケーションプログラムの入出力同期待ちによる処理遅延を回避することができる。

【 0 0 8 8 】 また、入出力制御装置にタイマ機能、入出力動作を終了させる手段、及び、入出力中断時のステータス報告手段(転送済みデータ量の報告) を設けるようにしたので、機能の実現が容易にできる。

【 0 0 8 9 】 また、複数のディスク装置で構成される論理ディスク装置において、この複数のディスク装置を並列動作させる場合に、ディスク装置間の入出力完了の同期を入出力動作の制限時間を設定することによりとるようにすることにより、ユーザの希望した時間までに、こ

の論理ディスク装置に対する入出力処理を完了することができる。つまり、論理ディスク装置を構成するディスク装置各々に対して、上位から与えられた同じ制限時間を設定して入出力命令を発行し、その完了を待ち、全てのディスク装置から応答があった時点で、各々の、入出力動作完了ステータスをまとめて上位に報告することにより、各々のディスク装置の入出力の完了状態はまちまちである可能性があるが、ユーザの希望した時間内に論理ディスク装置に対する入出力動作を終了させることができる。

【 0 0 9 0 】 さらに、ストライプ構成された論理ディスク装置に、複数のディスク装置間の同期方法として、この実施例で説明した同期方法を適用することにより、設定する制限時間により論理ディスク装置に対する入出力動作の応答時間は保証され、かつ、入出力動作の完了ステータスのばらつきの少ない、つまり、応答時間に見合ったデータ量設定であれば、全てのディスク装置における入出力動作が完了する可能性が高くなる。

【 0 0 9 1 】 実施例 4 . 本実施例は、本発明の入出力処理システムにおける高速処理を実現するためのさらに他の実施例である。本実施例では、上記実施例 1 ~ 実施例 3 で説明した入出力システムにおいて、入出力装置間のデータ転送を高速に行うための機構について説明する。

【 0 0 9 2 】 従来、入出力装置間でデータ転送する時には、データ転送制御を行うユーザプログラムデータのデータ空間に転送元からデータを読み込み、読み込んだデータを転送先の装置に書き込むといった操作が行われていた。ユーザプログラムの有するメモリ空間は、システムのメモリ管理下におかれ、データ転送の際システムのメモリ管理にとって負担となり、場合によってユーザプログラムのデータ空間にはデバイスから DMA 動作できないこともあり、これがデータ転送速度を結果的に低下させていた。これをシステムの物理メモリ上に、システムの間知しないデータ転送用のバッファメモリ領域を確保し、カーネル空間で実行される制御プロセスからコントロールして、データ転送することにより、メモリ管理へ負担を掛けることによる性能低下を防ぐことができる。

【 0 0 9 3 】 図 9 は、この高速転送を、論理ディスクを含まないデバイス間での実現するためのバスアービタ 1 0 5 の例で、図において 4 0 1、4 0 2 は各々データ転送の転送元、転送先となる高速デバイス 1 0 8 のローカルバス 1 0 7 上の I D が登録される、システム上のメモリ空間にマップされたバスアービタのソースレジスタ、デスティニレジスタ、4 0 3 はローカルバスを制御する論理回路である。図 1 0 は、論理ディスクを含まない入出力装置間の高速転送を行うデータ転送ドライバ 2 0 3 の処理の流れを示すフローチャート図である。図 1 1 は、論理ディスクを含む入出力装置間の高速転送を行う

バスアービタ105の処理の流れを示すフローチャート図である。

【0094】制御プログラム201により、データ転送ドライバ203に対し、転送元デバイス、その種別、及び、データオフセット、転送先デバイス、その種別、及び、データオフセット、そして、転送データ量を指定されるとデータ転送ドライバ203は図10のフローチャート図に示す論理に従って動作する。以後、図を参照しながら説明する。ステップ1001でデータ転送ドライバ203は、制御プログラム201より渡されたデバイス種別(通常デバイス、高速デバイス、論理ディスク装置のいずれか)から、転送元、あるいは、転送先のデバイスに論理ディスク装置を含んでいるか否かを判断する。論理ディスク装置を含む場合は、本実施例の場合、ボリュームマネージャ、及び、ディスクドライバに後述するバスアービタ105にインターフェースを設けていないので、バスアービタ105によるデータ転送機能は利用せず、データ転送ドライバ203が管理するシステムのメモリ102とデバイスドライバ(ディスクドライバ205、高速デバイスドライバ206)を利用して、

ステップ1004でデバイスを制御する。

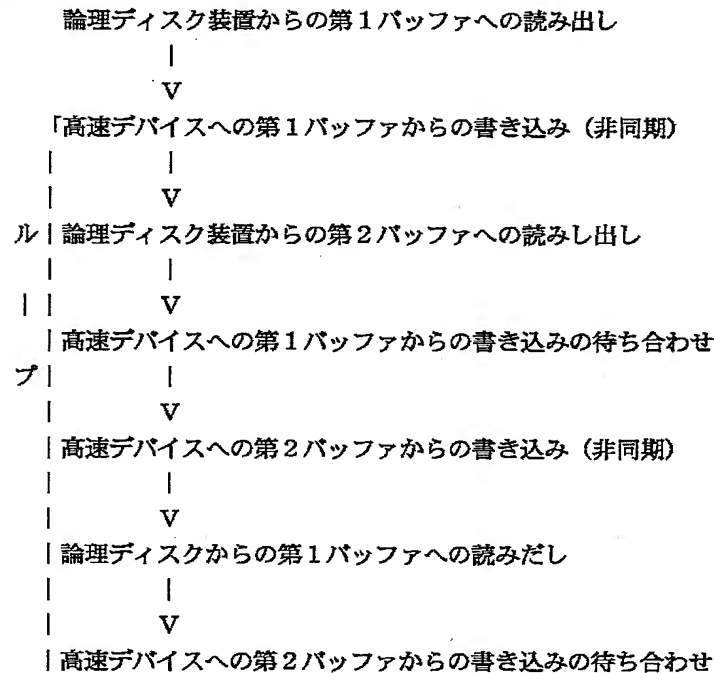
【0095】ステップ1004に示す処理につき説明す

る。ところで、データ転送ドライバ203は、データ転送処理のためのスタティックなメモリ領域をシステム上に確保している。このメモリはシステム初期化時に例えば2MB物理アドレス上連続に、ページング対象外ページとして確保されていて、カーネル空間として仮想空間にマップされている。これを2面バッファとして運用してデータ転送を行う。ストライプ構成された論理ディスク装置のチャンクサイズが128KBであれば8チャンク分の入出力サイズの指定が可能である。この処理の特徴的なところは制御プログラム201より発せられた、1本のコンテキストで非同期のダブルバッファ処理を行うことである。論理ディスク装置はボリュームマネージャ204により管理され、非同期インターフェースを持つ複数のディスク装置の入出力を管理するが、すでに説明したように、論理ディスク装置の入出力結果を同期的に利用者に返さねばならぬため非同期のデバイス同士の同期を取っている。このため、データ転送ドライバ203が論理ディスク装置にアクセスするときはコールして、その処理の終了を待たねばならない。従って、2面バッファの非同期運用を行うため本処理は、以下のようになっている。

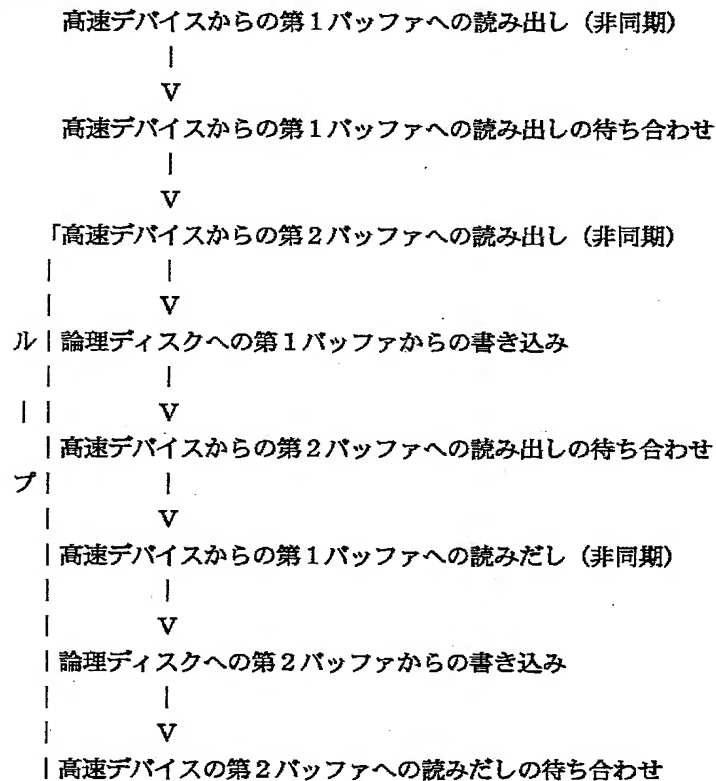
27

28

[転送元が論理ディスクの場合]



[転送元が高速デバイスの場合]



【0096】一方、データ転送対象デバイスに論理ディスク装置を含まない時、高速デバイス同士のデータ転送、あるいは、ディスクストライピングが有効でないケースのディスク装置と高速デバイスとのデータ転送がこれに当たるが、この場合は、データ転送機能を有するバ

スアービタ105を用いてデータ転送を行う。以下、この場合の動作について説明する。

【0097】ステップ1002において、データ転送ドライバ203によりデバイス間の転送速度の差がないと判断された場合、つまり、高速デバイス同士のデータ転

送の場合について説明する。高速デバイス108はDMAコントローラを有するので、バスアービタ105によりハンドシェイクすることでメモリ102を介在させずにデータ転送を行うことができる。ステップ1005でデータ転送ドライバはシステムのメモリ空間にマップされたバスアービタ105のソースレジスタ401、デステイニレジスタ402に、転送元、転送先となる高速デバイス108のローカルバス107上のID(デバイスアドレス)を登録する。さらに、転送元の高速デバイス108に対し指定された入出力サイズ、あらかじめ、システムから与えられた未使用の物理メモリ空間のアドレスでリード要求を発行する。また、転送先の高速デバイス108に対し、同様に、サイズ、アドレスを渡し、ライト要求を発行する。これにより、2つのデバイスからバスにバス獲得要求がだされるが、該IDはソースレジスタ401、デステイニレジスタ402に登録されている。さらに、ソースレジスタ401、デステイニレジスタ402にメモリ102の使用ビット(このビットについては後述する)は立っていないので、バスアービタ105は、ステップ1102、双方のバス要求が揃った時点で1104双方のバスグラント信号をアサートすることにより、メモリの介在なしのデータ転送が行われる。データ転送の終了の後、データ転送ドライバはアービタのソース/デステイニレジスタをクリアする。

【0098】データ転送ドライバ203は、ステップ1002でデバイス間の転送速度に差があると判断された場合、つまり、ディスク装置111と高速デバイス108の一方が転送元で他方が転送先となっている場合には、システムが提供するメモリアロケータにより、物理アドレス上連続でページアウトされないメモリ102のエリアを、入出力サイズ分確保し、ステップ1003でバスアービタ105のソースレジスタ401、デステイニレジスタ402に、各々のIDとメモリ使用ビットフラグをORして登録する。この後、ステップ1005で双方のデバイスに転送データサイズ、確保したメモリのアドレスを渡し、各々、リード/ライト要求を発行する。バスアービタ105は、ステップ1102でメモリ使用指示があるので、ソースデバイスからのバス要求に対しては、ステップ1103でバスグラント信号をアサートし、ソースデバイスのデータ転送終了を意味するビットをソースレジスタ401にORする。転送先のデバイスは、バス要求が認められるまでバス要求を繰り返す。バスアービタ105は、バスフリーでかつソースデバイスのデータ転送終了が転送終了ビットがONであることを確認すると、ステップ1106でグラント信号をアサートする。これにより、転送先デバイスは転送用に確保したメモリ102に転送されている転送元のデータを自デバイスに書き込む。これによりデータ転送が行われる。データ転送の終了の後、データ転送ドライバ203はバスアービタ105のソースレジスタ401、デ

ステイニレジスタ402をクリアし、データ転送のために確保したメモリ102を解放する。

【0099】以上のように、データ転送ドライバがデータ転送用の2面バッファを確保してデータ転送を行うようにしたので、入出力装置間のデータ転送を、より高速で行うことができる。即ち、システムのメモリ管理への負担の少ない、言い替えるとシステムのメモリ管理へ負担を掛けることによる性能低下を防ぐデータ転送方式を提供することができる。

【0100】また、各入出力装置の制御をする制御装置に設けられたDMAコントローラをコントロールメモリを介さずに、各々の入出力装置のデータ転送をバス上で同期させるコントローラを主たる入出力バスに備えることにより、無駄なメモリへのデータ転送を除くことができる。

【0101】また、データ転送を行う制御プログラム(データ転送ドライバ)が、入出力装置間の速度差を吸収するために、主たる入出力バスに接続されるデータ転送コントローラ(バスアービタ)にシステムからメモリを確保して与えるようにしたので、データ転送コントローラはこのメモリを利用してデバイス間の転送速度調節をすることができるようになり、入出力装置間のデータ転送を効率良く高速に行える。

【0102】実施例5. 本実施例は、実施例1～実施例4で構成された入出力システムにおいて、各デバイスの持つDMAコントローラを制御することにより(当然のことながら、デバイス間でのDMA転送レートは、一致させる必要はある)、メモリ102を介さずに各々のデバイスのデータ転送をバス上で同期させる制御ローカルバス107上に備えるようにして、無駄なメモリ102へのデータ転送を避ける例である。本実施例においては、各々のSCSIバス上の異なるディスク装置で4つのストライプディスクを持つ論理ディスク装置を構成することが多いが、データ整理のために論理ディスク装置の配置を変更する必要が発生することもある。この時、ボリュームマネージャ204は、データ移動先の対応するディスク装置が同一のSCSIバス上にあれば、全てのストライプディスク装置に対して移動先のディスク装置を複写先として、SCSIのコピーコマンドを発行することにより、データを複写し、複写元を未使用とすることにより、論理ディスク装置のデータ配置を変更することができる。

【0103】この実施例は、例えば、主たる入出力バス(ローカルバス)に複数のSCSIバスが接続され、各々のSCSIバスには複数のディスク装置が接続され、さらに、各々のSCSIバスから1つつつディスク装置を使用するように複数のストライプ構成の論理ディスク装置を構成したシステムに適用される。このようなシステムにおいて、論理ディスク装置間のデータ転送が必要となったときは、同一のバスに接続されたディスク装置

は、2つの論理ディスク装置は同一のストライプ幅で構成されているので、各々対応するストライプ同士で、SCSIのコピーコマンドを利用することにより、この論理ディスク装置間でデータ転送を行うことにより、高速で、システムへの負担の少ないデータ転送を実現できる。また、高速デバイスをSCSIコントローラ機能を備えた複数の入出力ポートを有する高速デバイスとすることによって、種類の異なる入出力装置間にも適用できる。

【0104】なお、上記各実施例においては、入出力インタフェースとしてSCSIを例にとり説明したが、インタフェースは特にSCSIである必要はなく、他のインタフェースを用いてもよいことは言うまでもない。また、上記実施例5は、ディスク装置間の複写をSCSIのコピーコマンドを用いて説明したが、ディスク制御装置にSCSIのコピーコマンド相当の機能を実現すればよい。

【0105】

【発明の効果】以上のように本発明によれば、システムに接続されているディスク装置の性能特性を性能収集手段がシステム管理者の与えるデータもしくは性能計測手段により収集し、収集したデータを基に、論理ディスク構成手段により論理ディスクを構成し、論理ディスク制御手段により論理ディスクを動作させるので、論理ディスクを構成する各ディスクの1ストライプ分のデータ入出力に要する応答時間が均等となり、論理ディスク装置の性能向上を図ることができる。

【0106】また、入出力命令の処理に要する入出力装置の実行時間を制限時間設定手段により規定し、制限時間が経過すると入出力命令を終了させる手段を設けたので、入出力動作の同期を時間を基準にコントロールすることが可能となり、アプリケーションプログラムでの余分な待ち時間を削減することができる、且つ大量のデータ転送を行う際の処理を効率良く行うことができる。

【0107】また、入出力制御装置に入出力命令の制限時間を設定するタイマとこのタイマが設定時間を経過すると実行中の入出力命令を終了させる手段を設けるようにしたので、処理の実現を容易に行うことができる。

【0108】また、複数のディスク装置で構成される論理ディスク装置において、この複数のディスク装置を並列動作させる場合に、ディスク装置間の入出力完了の同期を入出力動作の制限時間を設定することによりとるようにすることにより、ユーザの希望した時間までに、この論理ディスク装置に対する入出力処理を完了することができる。

【0109】さらに、ストライプ構成された論理ディスク装置に、複数のディスク装置間の同期方法として、この実施例で説明した同期方法を適用することにより、設定する制限時間により論理ディスク装置に対する入出力動作の応答時間は保証され、かつ、入出力動作の完了ス

テータスのばらつきの少ない、つまり、応答時間に見合ったデータ量設定であれば、全てのディスク装置における入出力動作が完了させることができる。

【0110】また、性能計測手段でディスク装置の性能を計測する際に、入出力命令に対する応答時間が最短あるいは最長となる条件を設定して行うようにしたので、論理ディスク構成装置が論理ディスクを構成する際に性能が最適になる論理ディスク装置を構成することができる。

【0111】また、システムにディスク装置を導入する際に、初期化に一部として性能計測手段により性能を計測するようにしたのでシステムとしての効率を向上させることができる。

【0112】また、性能計測手段に入出力装置の性能だけでなく、入出力バスの性能を計測する手段を設け、性能収集手段を構成するようにしたのでシステム全体の入出力性能を向上させることができる。

【0113】また、論理ディスク構成手段が、1台のディスク装置に対する1回の入出力命令で行うデータ量をそのディスク装置の1トラックに収まるように論理ディスク構成データを構成するようにしたので、論理ディスク装置を構成するディスク装置のシーク方向が順方向と逆方向での性能差がなくなる。

【0114】また、論理ディスク構成手段が内外周での性能に差があるディスク装置が論理ディスク装置に含まれているときは、内外周を交互に組み合わせて論理ディスクを構成するようにするので、各ストライプに対して均等な処理時間となり、処理効率を向上させることができる。

【0115】また、論理ディスク制御手段がディスク装置に指令を出すときに、各ディスク装置のヘッド位置を考慮して入出力実行時間が最少となるように制御するので処理効率が向上する。

【0116】また、同一のディスク装置に対する複数の入出力命令を纏めるようにしたので、命令実行に伴う処理時間を減少させることができるので入出力処理を向上させることができる。

【0117】また、論理ディスク装置の構成を自動的に行うようにしたので、論理ディスクの構成が容易になる。

【0118】また、データ転送ドライバが入出力命令の対象とする入出力装置に論理ディスク装置が含まれている場合は、2面バッファをシステムメモリ上に確保してデータ転送を行うようにしたので、アプリケーションプログラムでメモリを確保する必要がなく、通常アプリケーション上のバッファメモリを介して行われていたデバイス間のデータ転送を、システムの仮想記憶に負担を掛けないページ管理外のメモリを使用して、あるいは、若干の特殊ハードウェアを付加することによりメモリを使用せずにデータ転送することを可能にし、大量高速デー

10

20

30

40

50

タ転送をシステムへの負担を軽減し、なおかつ、高速に実行できる。

【0119】また、入出力バスの制御をするバスアービタに入出力装置の並列動作をさせる手段を設け、入出力命令の対象に論理ディスク装置が含まれていない場合には、このバスアービタを動作させるようにしてシステムメモリを使用せずにデータ転送をすることができるようにしたので、データ転送を高速にできる。

【0120】また、データ転送ドライバに入出力命令の対象となる入出力装置間に速度性能差がある場合には、システムメモリにデータ転送用のバッファを確保してバスアービタを動作させるようにしたので、アプリケーションでメモリを確保する必要のない効率のよいデータ転送を行うことができる。

【0121】また、ディスク装置を制御するディスク制御装置に同一バスに接続されたディスク装置間の複写手段を設けるようにしたので、論理ディスク装置のバックアップが容易に実現できる。

【図面の簡単な説明】

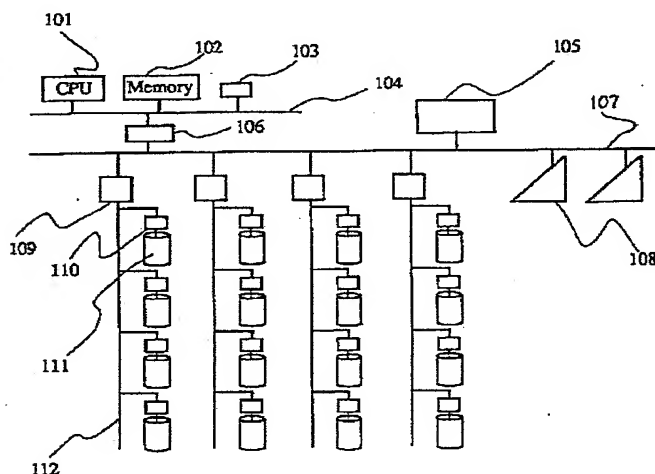
【図1】 本発明による入出力処理システムの適用されるシステムのハードウェアの構成の一例を示す図である。

【図2】 本発明による入出力処理システムの適用されるシステムのソフトウェアの構成の一例を示す図である。

【図3】 本発明による性能計測手段によって作成されたディスクの性能データの一例を示す図である。

【図4】 本発明によるボリュームマネージャの行う論理ディスク装置を自動的に構成するための処理の流れを

【 図1 】



示す図である。

【図5】 本発明による論理ディスク管理データの一例を示す図である。

【図6】 本発明による論理ディスク構成データの一例を示す図である。

【図7】 本発明におけるディスクコントローラの構成を示すブロック図である。

【図8】 本発明における論理ディスク装置の動作状況の一例を示す図である。

【図9】 本発明におけるバスアービタの構成を示すブロック図である。

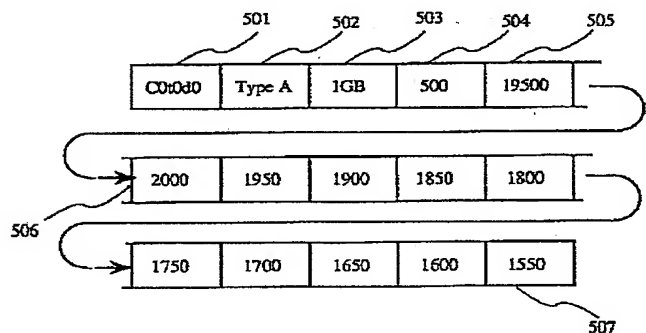
【図10】 本発明によるデータ転送ドライバの処理の流れの一例を示す図である。

【図11】 本発明におけるバスアービタの処理の流れの一例を示す図である。

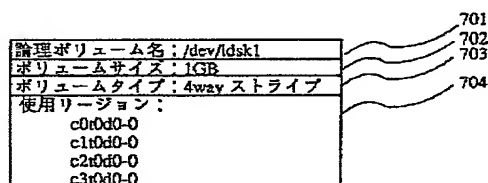
【符号の説明】

105 バスアービタ、108 高速デバイス、109 SCSI ホストアダプタ、110 ディスクコントローラ、201 データ転送制御プログラム、203 データ転送ドライバ、204 ボリュームマネージャ、208 性能計測手段、209 性能データ、210 ディスク管理データ、211 論理ディスク構成データ、212 論理ディスク入出力インターフェース、213 論理ディスク構成手段、302 SCSI コントロールロジック、304 ディスクコントロールロジック、305 タイマ、401 ソースレジスタ、402 デスティニレジスタ、403 バスアービトリションロジック。

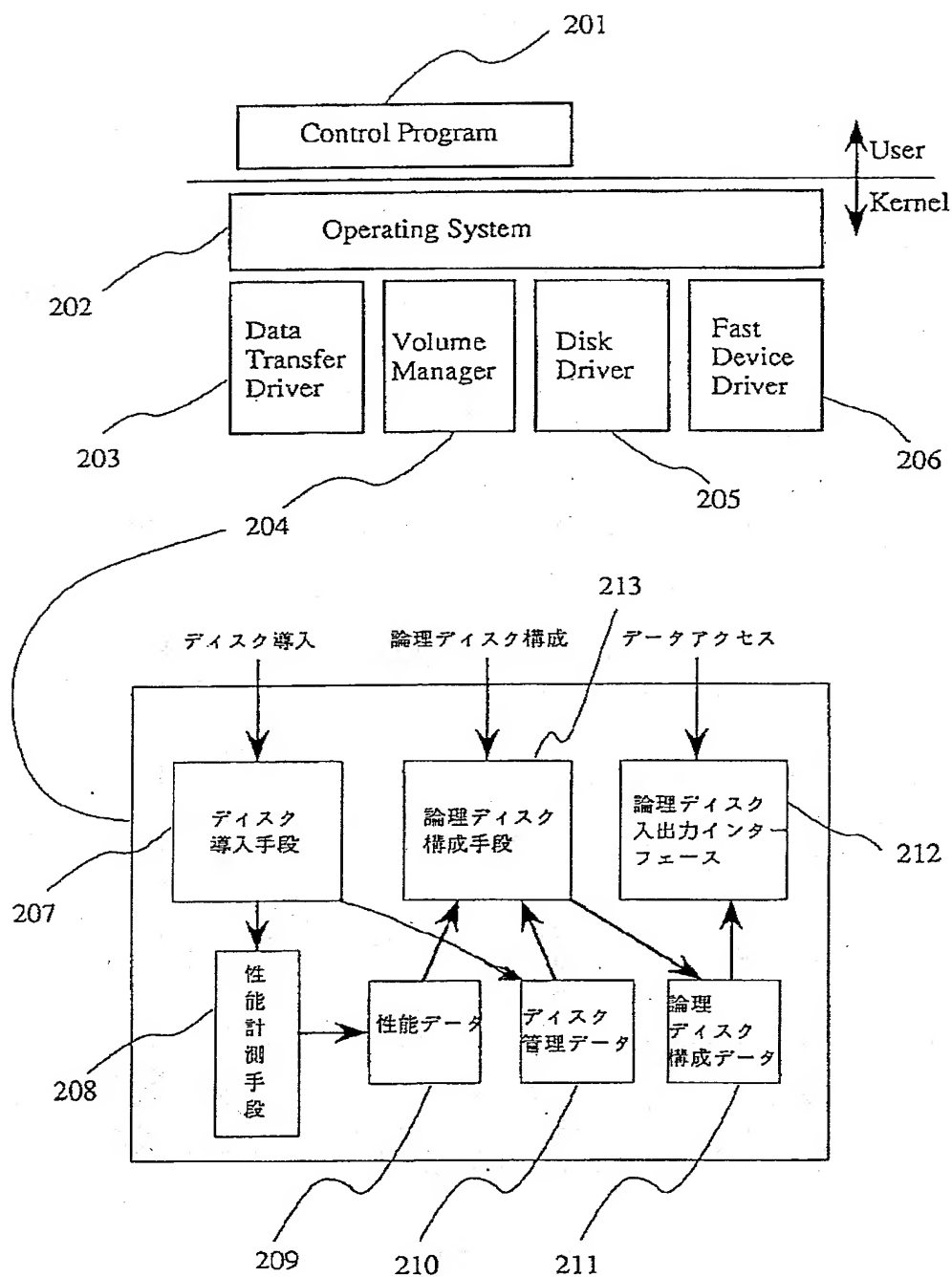
【 図3 】



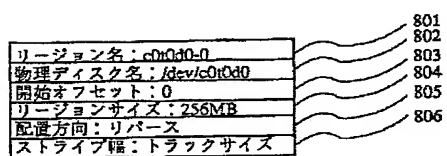
【 図5 】



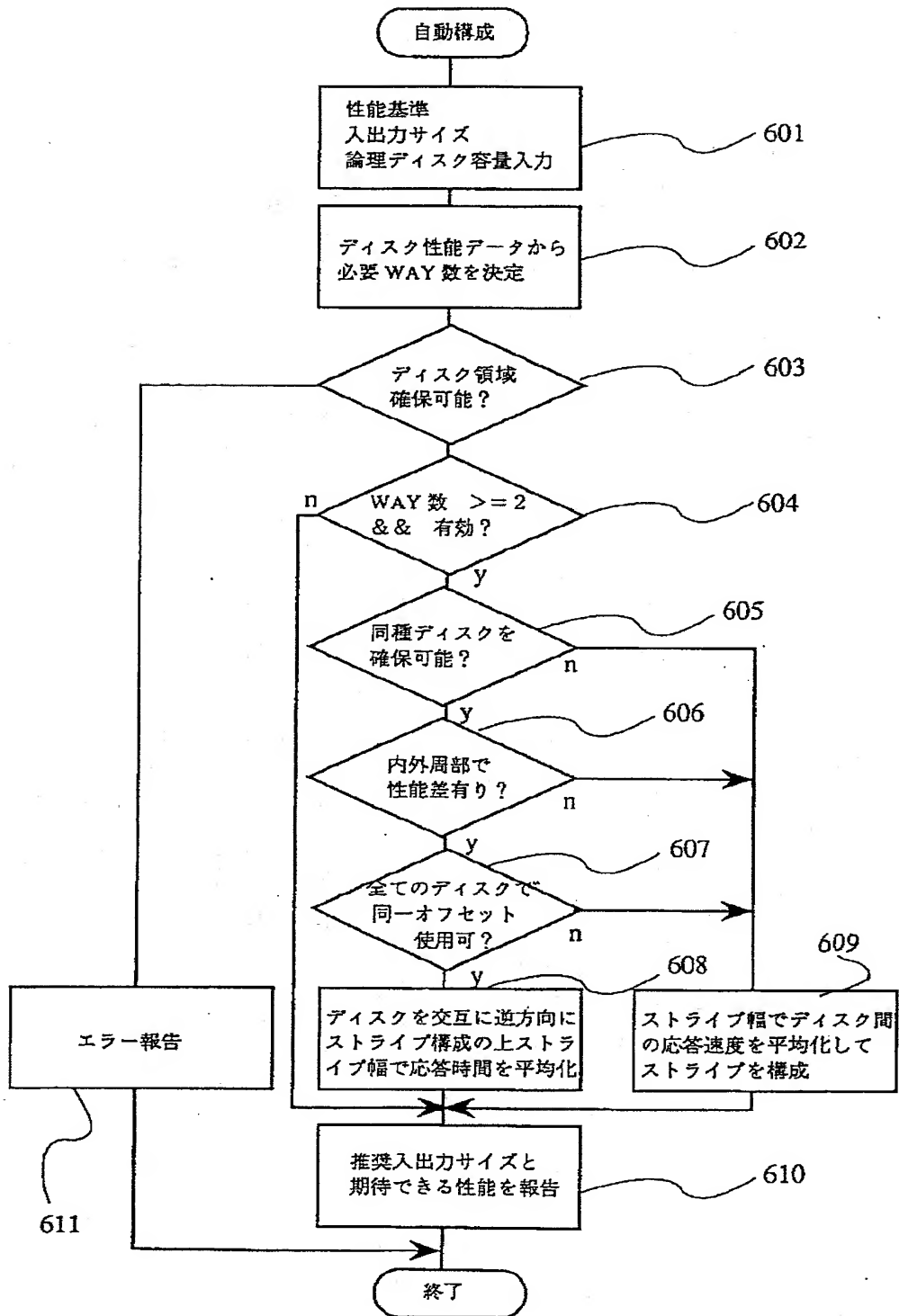
【 図2 】



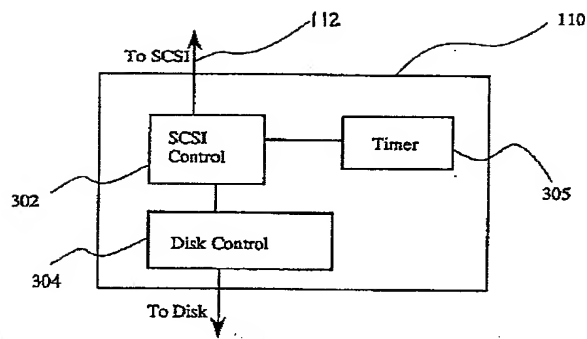
【 図6 】



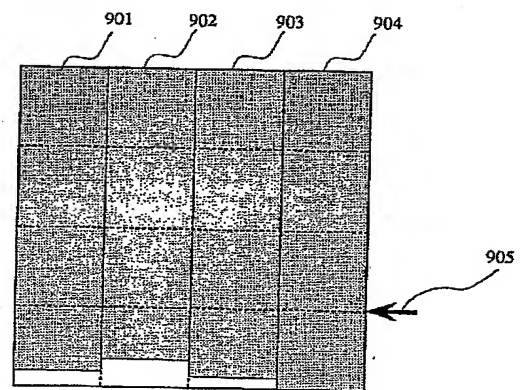
【 図4 】



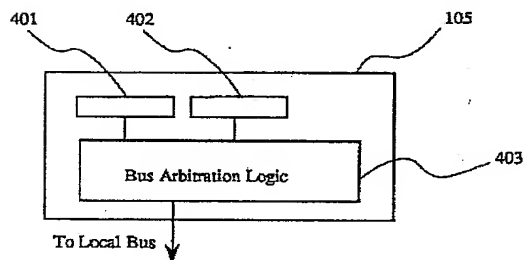
【 図7 】



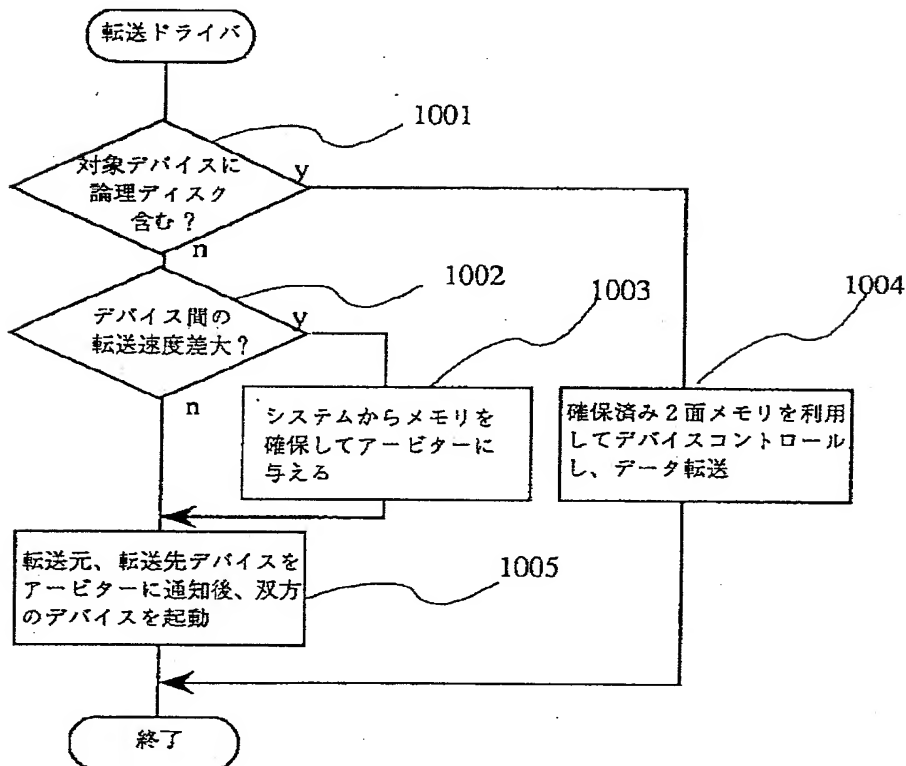
【 図8 】



【 図9 】



【 図10 】



【 図1 1 】

